



Artificial Intelligence in Data Processing

Alexander Measure

Economist

Office of Safety and Health Statistics

FedCASIC 2013





Survey of Occupational Injuries and Illnesses (SOII)

- Annual Survey
- 250,000 cases
- Text narratives for:
 - ▶ Occupation
 - ▶ Injury/Illness Characteristics

Example Case

Job title: janitor

Occup: 37-2011 (Janitor)

What was the employee doing just before the incident?

mopping floor in gym

Nature: 111 (Fracture)

What happened?

slipped on wet floor and fell

Event: 422 (Fall, slipping)

What part of the body was affected?

fractured right arm

Part: 420 (Arm)

What object directly harmed the employee?

wet floor

Source: 6620 (Floor)

Limitations of Human Coding

- Time consuming
- Very difficult
 - ▶ Detailed classifications
 - ▶ Ambiguous data
- Inconsistent coding



Artificial Intelligence

- Text Classification
 - ▶ Rules
 - ▶ Machine Learning
- Ex. Census
 - ▶ Rule based system: 192 person-months
 - ▶ Machine learning: 4 person-months
 - Source: Creecy et al., 1992



Machine Learning

■ 3 Steps

1. Select feature representation
2. Select model
3. Fit model to data

Feature Representation

- Each feature corresponds to a word
- Job title: “assistant nurse”
 - ▶ $X_{\text{nurse}} = 1$
 - ▶ $X_{\text{assistant}} = 1$
 - ▶ $X_{\text{mechanic}} = 0$

Logistic Regression Model

$$P(\text{code} = A) = \frac{\exp(w_{a1}x_{\text{nurse}} + w_{a2}x_{\text{assistant}} + w_{a3}x_{\text{mechanic}})}{Z}$$

$$P(\text{code} = B) = \frac{\exp(w_{b1}x_{\text{nurse}} + w_{b2}x_{\text{assistant}} + w_{b3}x_{\text{mechanic}})}{Z}$$

$$P(\text{code} = C) = \frac{\exp(w_{c1}x_{\text{nurse}} + w_{c2}x_{\text{assistant}} + w_{c3}x_{\text{mechanic}})}{Z}$$

Autocode

- New Case: "assistant mechanic"

- ▶ $X_{\text{nurse}} = 0$

- ▶ $X_{\text{assistant}} = 1$

- ▶ $X_{\text{mechanic}} = 1$

- $P(\text{code}=A) = .05$

- $P(\text{code}=B) = .25$

- $P(\text{code}=C) = .70$ ←

Does it work?

- Train on 250k
- Test on 10k

	Computer to Original (%)
Occupation	80
Part	81
Nature	80
Event	49
Source	60

Humans vs. Computer

- 1000 cases
- 3 re-coders

	Human to Original (%)	Computer to Original (%)
Occupation	66	80
Part	82	81
Nature	76	80
Event	47	49
Source	56	60



Limitations

- Training data
- New words
- Linguistic complexity



Applications

- Review
- Assisted Coding
- Partial Autocoding



Useful Resources

- Free Machine Learning Class
 - ▶ <https://www.coursera.org/course/ml>
- Free Software
 - ▶ Python (Scikit-Learn)
- Free Support
 - ▶ MetaOptimize Q&A
 - ▶ Stack Overflow

Contact Information

Alexander Measure

Economist

Office of Safety and Health Statistics

[*www.bls.gov/iif*](http://www.bls.gov/iif)

202-691-6185

measure.alex@bls.gov

